

# UCLA

## UCLA Previously Published Works

### Title

Global circulation patterns of seasonal influenza viruses vary with antigenic drift.

### Permalink

<https://escholarship.org/uc/item/2rw1g49w>

### Journal

Nature, 523(7559)

### ISSN

0028-0836

### Authors

Bedford, Trevor  
Riley, Steven  
Barr, Ian G  
et al.

### Publication Date

2015-07-01

### DOI

10.1038/nature14460

Peer reviewed

Published in final edited form as:

*Nature*. 2015 July 9; 523(7559): 217–220. doi:10.1038/nature14460.

## Global circulation patterns of seasonal influenza viruses vary with antigenic drift

Trevor Bedford<sup>1</sup>, Steven Riley<sup>2,3</sup>, Ian G. Barr<sup>4</sup>, Shobha Broor<sup>5</sup>, Mandeep Chadha<sup>6</sup>, Nancy J. Cox<sup>7</sup>, Rodney S. Daniels<sup>8</sup>, C. Palani Gunasekaran<sup>9</sup>, Aeron C. Hurt<sup>4,10</sup>, Anne Kelso<sup>4</sup>, Alexander Klimov<sup>7</sup>, Nicola S. Lewis<sup>11</sup>, Xiyang Li<sup>12</sup>, John W. McCauley<sup>8</sup>, Takato Odagiri<sup>13</sup>, Varsha Potdar<sup>6</sup>, Andrew Rambaut<sup>3,14,15</sup>, Yuelong Shu<sup>12</sup>, Eugene Skepner<sup>11</sup>, Derek J. Smith<sup>11,16</sup>, Marc A. Suchard<sup>17,18,19</sup>, Masato Tashiro<sup>13</sup>, Dayan Wang<sup>12</sup>, Xiyang Xu<sup>7</sup>, Philippe Lemey<sup>20</sup>, and Colin A. Russell<sup>21</sup>

<sup>1</sup>Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA <sup>2</sup>MRC Centre for Outbreak Analysis and Modelling, Department of Infectious Disease Epidemiology, School of Public Health, Imperial College London, London, UK <sup>3</sup>Fogarty International Center, National Institutes of Health, Bethesda, MD, USA <sup>4</sup>World Health Organization (WHO) Collaborating Centre for Reference and Research on Influenza, Melbourne, Australia <sup>5</sup>SGT Medical College, Hospital and Research Institute, Village Budhera, District Gurgaon, Haryana, India <sup>6</sup>National Institute of Virology, Pune, India <sup>7</sup>WHO Collaborating Center for Reference and Research on Influenza, Centers for Disease Control and Prevention, Atlanta, GA, USA <sup>8</sup>WHO Collaborating Center for Reference and Research on Influenza, Medical Research Council National Institute for Medical Research (NIMR), London, UK <sup>9</sup>King Institute of Preventive Medicine and Research, Guindy, Chennai, India <sup>10</sup>Melbourne School of Population and Global Health, University of Melbourne, Parkville VIC 3010, Australia <sup>11</sup>Department of Zoology, University of Cambridge, Cambridge, UK <sup>12</sup>WHO Collaborating Center for Reference and Research on Influenza, National Institute for Viral Disease Control and Prevention, China CDC, Beijing, China <sup>13</sup>WHO Collaborating Center for Reference and Research on Influenza, National Institute of Infectious Diseases, Tokyo, Japan <sup>14</sup>Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK <sup>15</sup>Centre for Immunology, Infection and Evolution, University of Edinburgh, Edinburgh, UK <sup>16</sup>Department of Viroscience, Erasmus Medical Center, Rotterdam, The Netherlands <sup>17</sup>Department of Biostatistics, UCLA Fielding School of Public Health, University of California, Los Angeles, CA, USA <sup>18</sup>Department of Biomathematics David

Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: [http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Correspondence and requests for materials should be addressed to C.A.R. (car44@cam.ac.uk).

**Author Contributions** C.A.R. and T.B. conceived the research. C.A.R. and T.B. drafted the manuscript with substantial support from P.L. and S.R. I.G.B., S.B., M.C., N.J.C., R.S.D., C.P.G., A.C.H., A.K., A.K., X.L., J.W.M., T.O., V.P., Y.S., M.T., D.W., and X.X. coordinated and produced the influenza surveillance data. T.B. performed the modeling and data analyses along with C.A.R., S.R., P.L., M.A.S. and A.R.. T.B. created the figures. All authors discussed the results and contributed to the revision of the final manuscript.

The authors declare no competing financial interests.

**Online Content** Methods, along with Extended Data Tables and Figures, and Supplementary Information are available in the online version of the paper; References unique to these sections appear only in the online paper. Accession numbers of all sequences used in this study are available to download as Extended Data.

Geffen School of Medicine at UCLA, University of California, Los Angeles, CA, USA

<sup>19</sup>Department of Human Genetics, David Geffen School of Medicine at UCLA, University of California, Los Angeles, CA, USA <sup>20</sup>Department of Microbiology and Immunology, Rega Institute, KU Leuven – University of Leuven, Leuven, Belgium <sup>21</sup>Department of Veterinary Medicine, University of Cambridge, Cambridge, UK

## Abstract

Understanding the spatio-temporal patterns of emergence and circulation of new human seasonal influenza virus variants is a key scientific and public health challenge. The global circulation patterns of influenza A/H3N2 viruses are well-characterized<sup>1-7</sup> but the patterns of A/H1N1 and B viruses have remained largely unexplored. Here, based on analyses of 9,604 hemagglutinin sequences of human seasonal influenza viruses from 2000–2012, we show that the global circulation patterns of A/H1N1 (up to 2009), B/Victoria, and B/Yamagata viruses differ substantially from those of A/H3N2 viruses. While genetic variants of A/H3N2 viruses did not persist locally between epidemics and were reseeded from East and Southeast (E-SE) Asia, genetic variants of A/H1N1 and B viruses persisted across multiple seasons and exhibited complex global dynamics with E-SE Asia playing a limited role in disseminating new variants. The less frequent global movement of influenza A/H1N1 and B viruses coincided with slower rates of antigenic evolution, lower ages of infection, and smaller less frequent epidemics compared to A/H3N2 viruses. Detailed epidemic models support differences in age of infection, combined with the less frequent travel of children, as likely drivers of the differences in the patterns of global circulation, suggesting a complex interaction between virus evolution, epidemiology and human behavior.

---

Owing to the frequency and severity of human seasonal influenza A H3N2 virus epidemics, recent work has focused on the global circulation dynamics of H3N2 viruses<sup>1-7</sup>. Studies have shown that, each year, H3N2 epidemics worldwide result from the introduction of new genetic variants from East and Southeast (E-SE) Asia, where viruses circulate via a network of temporally overlapping epidemics<sup>1,2,4,5</sup>, rather than local persistence<sup>1,3,6,7</sup>. In addition to H3N2, H1N1 viruses and two antigenically diverged lineages of influenza B viruses, B/Victoria/2/1987-like (Vic) and B/Yamagata/16/1988-like (Yam), circulate among humans with lower but substantial disease burdens<sup>8,9</sup>. Despite their importance, the global circulation dynamics of former seasonal H1N1 viruses (preceeding the 2009 pandemic) and B viruses have been largely neglected.

Given that influenza A and B viruses cause similar symptoms and evolve by similar mechanisms of immune escape, we hypothesized that each would follow similar patterns of global circulation, with new genetic variants originating in E-SE Asia that rapidly replace existing genetic variants. To test this hypothesis we compared the global circulation patterns of the hemagglutinin (HA) genes of H3N2, former seasonal H1N1, Vic, and Yam viruses. We assembled datasets of HA sequences with complete HA1 domains for each subtype from the World Health Organization Global Influenza Surveillance and Response System and the Influenza Research Database<sup>10</sup> covering 2000–2012. To reduce the impact of surveillance

biases, we subsampled these data to more equitable spatiotemporal distributions, resulting in datasets comprising 4006 H3N2, 2144 H1N1, 1999 Vic, and 1455 Yam HA sequences (Extended Data Fig. 1). Though deficient in viruses from Africa and Eastern Europe, these are the most geographically and temporally comprehensive seasonal influenza virus datasets assembled to date.

By estimating temporally-resolved phylogenetic trees for each subtype, we revealed faster rates of nucleotide mutation and amino acid substitution in H3N2 and H1N1 than in the B viruses (consistent with previous work<sup>11,12</sup>), but more genealogical diversity in B viruses than A viruses (Extended Data Table 1). This inverse relationship between evolutionary rate and genealogical diversity is expected if increased mutation rate correlates with antigenic drift<sup>13</sup> and drives increased adaptive evolution, thus purging HA genetic diversity<sup>14</sup>. By inferring geographic ancestry using Bayesian phylogeographic methods<sup>15</sup>, we found a consistent pattern for H3N2 viruses (Fig. 1a) in which viruses worldwide rapidly coalesce to the trunk of the tree (average time to trunk = 1.42 years), with trunk viruses mostly originating from E-SE Asia (Extended Data Fig. 2a). This finding is consistent with previously reported patterns<sup>1,2,4,5</sup>, with E-SE Asia acting as the source population for epidemics worldwide.

In addition to China and Southeast Asia, India frequently contributed viruses to the trunk of the tree suggesting that the global circulation of H3N2 viruses is maintained by an E-SE Asian network that includes India. India's role in the global dissemination of H3N2 viruses may have been similar historically, but India-wide influenza surveillance only began in 2004. There were brief periods, notably the 2007–2008 Northern Hemisphere winter, when regions outside E-SE Asia contributed to the trunk of the H3N2 tree. However, these instances were rare and trunk viruses from outside E-SE Asia descended directly from viruses within E-SE Asia (Fig. 1a). Quantifying the average ancestry of strains from each geographic region in the 3 years prior to sampling showed prominent roles for China, India, and Southeast Asia in seeding epidemics in all regions (Extended Data Fig. 3).

Surprisingly, the global circulation patterns of former seasonal H1N1 viruses differed substantially from those observed for H3N2 viruses (Fig. 1). Like H3N2, most lineages of H1N1 viruses eventually coalesced with viruses from E-SE Asia and India. However, this coalescence was slower than for H3N2 viruses with prolonged co-circulation of geographically segregated H1N1 lineages (Fig. 1b, Extended Data Figs. 3 and 4). Geographic segregation of H1N1 viruses was particularly pronounced beginning in 2004/2005, with the emergence of three co-circulating genetic lineages (Fig. 1b, nodes 1-3) that each independently acquired HA mutations leading to antigenic evolution from the A/New Caledonia/20/1999-like phenotype to the A/Solomon Islands/3/2006-like phenotype. These lineages circulated in Southeast Asia (node 1), China (node 2) and India (node 3), with the Indian lineage eventually spreading worldwide prior to the emergence of H1N1pdm09 viruses.

Phylogeographic analyses of B Vic and Yam viruses revealed further differences from H3N2 viruses with lineages frequently circulating outside of E-SE Asia for several years without evidence of seeding from E-SE Asia (Fig. 1c,d). Prominent examples include the

seeding of the North American 2006/2007 Vic season directly from 2005/2006 North American viruses and the seeding of the North American 2001/2002 Yam season directly from 2000/2001 North American viruses (Extended Data Fig. 4). Similarly, lineages of viruses within E-SE Asia commonly circulated exclusively in E-SE Asia for more than 1 year. These long circulating E-SE Asian lineages were most apparent for Vic viruses where two lineages (Fig. 1c, nodes 1 and 2) persisted independently in China and SE Asia for over 5 years without spreading to other regions and led to the co-circulation of three distinct Vic antigenic variants in different parts of the world during 2007/2008 (Extended Data Fig. 5a).

Patterns of persistence of genetic variants differed by subtype and region, with H3N2 viruses persisting regionally for an average of ~6 months, H1N1 for ~9 months, Vic for ~13 months and Yam for ~12 months. H3N2 viruses showed comparably short durations of persistence across the world (Fig. 1), with the exceptions of India and China. Patterns within China were characterized by North and South lineages contributing jointly to persistence as combining North and South phylogeny nodes resulted in substantially greater persistence estimates than from North or South lineages alone (Fig. 1). For H3N2, evidence for joint contributions to persistence by region pairs that exclude China is comparatively weak (Extended Data Fig. 6a, Supplementary Information). For Vic and Yam, the mean duration of persistence was longer than for H3N2 or H1N1 in most regions, particularly in India and China where mean durations were >2 years (Fig. 1, Extended Data Fig. 4). Duration of regional persistence correlated with the proportion of virus originating from that region (Extended Data Fig. 6b) and observed phylogeographic patterns were robust to subsampling assumptions (Supplementary Information, Extended Data Table 2).

To investigate differences in the global migration patterns of H3N2, H1N1 and B viruses, we used the spatiotemporally-resolved phylogenies to estimate the amounts of virus movement between regions (Fig. 2). Rates of movement between pairs of regions were highly correlated between viruses with Spearman correlation coefficients ranging from 0.65 (H3N2 vs Yam) to 0.75 (H3N2 vs H1N1), suggesting similar global connectivity networks for all viruses. However, while the overall structure of the migration network was similar, H3N2 viruses moved between regions more frequently than H1N1 and B viruses (migration events per lineage per year H3N2 = 1.96, H1N1 = 1.27, Vic = 0.93, Yam = 0.97, Extended Data Table 1).

We hypothesized a relationship between rates of global movement and rates of antigenic drift: though rates of genetic evolution were similar for H3N2 and H1N1 viruses, both H1N1 and B viruses evolved antigenically more slowly than H3N2 viruses<sup>13</sup> (Extended Data Table 1). We also hypothesized that lower rates of immune escape for B and H1N1 compared with H3N2 would lead to: younger average ages of infection as children increasingly comprise the largest pool of susceptible individuals; and smaller, less frequent epidemics owing to smaller populations of susceptible individuals<sup>13</sup>. These differences are consistent with results from several community-based cohort studies that found that children were more frequently infected with B viruses than were adults<sup>8,16,17</sup>. Age of infection data covering 2002–2011 from Australia show that H1N1 and B viruses infect younger individuals than H3N2 viruses (Extended Data Fig. 5b-d, median age of infection H3N2 = 30y, H1N1 = 20y,

B = 16y) and epidemiological data from Australia and the United States show reduced size and frequency of H1N1 and B epidemics compared to H3N2 (Extended Data Fig. 5f-i).

Differences in age of infection may explain differences in global circulation as children travel long distances much less frequently than adults (Extended Data Fig. 5e). A previous study hypothesized that age-specific patterns of infection could lead to differences in contact rates and the spread of influenza types within the United States over the course of a single season<sup>18</sup>. Here, we hypothesized that differential global air travel provides a plausible mechanism by which H1N1 and B viruses show increased genetic differentiation and reduced rates of global migration across multiple seasons, compared to H3N2 viruses.

To test the impact of differences in age distribution of infection on global patterns of virus movement, we constructed a multi-patch transmission model. We modeled two scenarios for host movement: 1. age-independent mixing between patches; 2. age-stratified mixing with host movement derived from air travel passenger age data (Extended Data Fig. 5e). In the age-independent scenario, model parameters only differed in rate of antigenic mutation, leading to differences in observed rates of antigenic drift among viruses and hence epidemic size and frequency (Extended Data Fig. 7). Faster antigenic drift resulted in greater incidence and more adult infections (Fig. 3a,b), but only modest differences in virus lineage movement (Fig. 3c, B-like viruses differ from H3-like viruses by a factor of 1.2), consistent with slightly faster spread of antigenically novel strains. However, age-stratified mixing between patches intensified the effect of antigenic drift on migration rate and created differences in rates of movement between patches more consistent with those observed for H3N2 vs H1N1 and B (Fig. 3c, B-like viruses differ from H3-like by a factor of 1.6). In the scenario with faster antigenic drift, infections were more mobile due to greater frequency of adult infection, causing a knock-on effect on rates of viral movement. The model also suggests that the differences in patterns of regional persistence observed in the phylogenies might be shaped by a combination of differences in rates of antigenic evolution and variation in amplitude of epidemic seasonality, with slowly evolving viruses persisting longer than rapidly evolving viruses at low amplitudes of seasonal forcing (Extended Data Fig. 8a, Supplementary Information).

In the model, varying transmission rate rather than antigenic mutation rate also resulted in differences in the observed rate of antigenic drift, with higher transmission resulting in faster drift (Extended Data Fig. 8b). The relationship between antigenic drift rate and migration rate is similar regardless of whether drift is modulated by mutation rate or transmission rate (Extended Data Fig. 8b). This finding is in line with theoretical work showing that epidemiological processes can influence influenza virus evolution<sup>19,20</sup>. However, there are important virological differences between influenza viruses that are likely to impact the efficiency and tempo at which antigenic variation is generated and fixed, which could in turn affect epidemiology<sup>21-24</sup> (Supplementary Information).

Regardless of the underlying drivers, there is a remarkable correspondence in model behavior, quantified as a stable relationship between observable rate of antigenic drift and global circulation patterns. The patterns of epidemic spread observed here suggest that

differences in ages of infection could explain patterns of global circulation across a variety of human viruses.

## Methods

### Sequence data

Hemagglutinin (HA) coding sequences for influenza A H3N2 viruses, former seasonal H1N1 viruses (preceding the 2009 pandemic), and influenza B virus lineages Victoria (Vic) and Yamagata (Yam) collected by the World Health Organization (WHO) Global Influenza Surveillance and Response Network including the National Institute of Virology, Pune, India between 2000 and 2012 were combined with human seasonal influenza virus sequences (minimum length = 984 base pairs) covering 2000 to 2012 from the Influenza Research Database<sup>10</sup>. After removing duplicate strains and strains overly divergent based on root-to-tip distances, the data set contained 9139 H3N2 sequences, 3789 H1N1 sequences, 2577 Vic sequences and 1821 Yam sequences. Sampling locations for these sequences were parsed from strain names. Sequences were grouped into 9 geographic regions: USA/Canada, South America, Europe, India, North China, South China, Japan/Korea, Southeast Asia and Oceania. Specifics of this partitioning are shown in Extended Data Figure 1. Groups were chosen to maximize available sequences within each region while still providing enough geographic diversity to ensure nearly global coverage. Sequences from Africa, Central America, the Middle East and Russia were excluded because of a lack of sufficient numbers of sequences to provide comparable estimates to other regions.

In the raw sequence data, some regions, such as the USA, were over-represented. Additionally, more recent years were over-represented compared to years at the start of the study period. In order to control for these sampling biases, we subsampled the raw data randomly by location and time to create a more equitable spatiotemporal distribution. The USA had consistently more sequences available every year from 2000 to 2012, thus in order to maintain similar total numbers of sequences for each region across the entire study period it was necessary to sample fewer sequences per year from the USA. We selected 50 sequences per region per year (40 for USA/Canada) for H3N2 and 80 sequences per region per year (45 for USA/Canada) for H1N1, Vic and Yam. This subsampling resulted in largely similar sequence counts across years and across regions for each virus, but overall more H3N2 sequences than H1N1 or B sequences, with 4006 H3N2 sequences, 2144 H1N1 sequences, 1999 Vic sequences and 1455 Yam sequences (Extended Data Fig. 1). When selecting subsampled sequences we first selected sequences with full day-month-year collection dates and then longer sequences over sequences with less precise dates or shorter sequences. HA sequence data for 1630 H3N2 isolates, 1600 H1N1 isolates, 1394 Vic isolates and 881 Yam isolates have been deposited in the Influenza Research Database<sup>10</sup> and accession numbers for all sequences used provided as Supplementary Information.

### Phylogeographic inference

Time-resolved phylogenetic trees were estimated for H3N2, H1N1, Vic and Yam using BEAST v1.8.1<sup>25</sup> and incorporated the SRD06 nucleotide substitution model<sup>26</sup>, a coalescent demographic model with constant effective population size and a strict molecular clock



across branches. A strict molecular clock was chosen based on finding strong correlations between date of sampling and evolutionary distance in all datasets, as estimated by Path-O-Gen v1.4<sup>27</sup>. Using a strict clock also reduced the risk of model over-parameterization (for example, for the complete H3N2 data set with a relaxed clock, there would be  $2 \times 4006 - 2 = 8010$  branch-specific rates). Samples with imprecise dates (known only to the month or to the year) had their dates of sampling estimated assuming a uniform prior within the known temporal bounds<sup>28</sup>. Markov Chain Monte Carlo (MCMC) was run for 600 million steps and trees were sampled every 5 million steps after allowing a burn-in of 100 million steps, yielding a total sample of 100 trees for H1N1, Vic and Yam. With significantly more samples, H3N2 required a longer chain to converge. Here, MCMC was run in parallel for 2 chains, each with 650 million steps sampled every 3 million steps with a burn-in of 500 million steps and samples across chains combined, yielding a total of 100 sampled trees. These trees were treated as independent draws from the posterior space of trees when subsequently used in the robust counting and phylogeographic analyses<sup>29</sup>. Evolutionary rates in Extended Data Table 1 were estimated using the ‘renaissance’ counting methods of Lemey et al.<sup>30</sup>.

Phylogeographic patterns were estimated using a discrete-state continuous time Markov chain (CTMC) model, in which transition rates were estimated between each pair of regions<sup>15</sup>. We assumed a non-reversible transition model<sup>31</sup> consisting of 72 separate rate parameters, each with a Bayesian stochastic search variable selection (BSSVS) indicator variable, and a separate overall rate of geographic transition. We assumed an exponential prior with mean of 1 for each transition rate, a negative binomial prior with mean of 9 and standard deviation of 9 for the total number of non-zero rates and an exponential prior with mean of 1 migration event per lineage per year for the overall geographic transition rate. MCMC was run for 12 million steps with a burn-in of 2 million steps, and parameters sampled every 10,000 steps and trees sampled every 100,000 steps, yielding a total sample of 1000 parameter states and 100 trees on which estimates were based. Pairwise migration rate estimates had an effective sample size (ESS) of 350 at the minimum and most had ESS greater than 500.

This procedure yielded posterior trees with the geographic states of internal nodes resolved. We analyzed these posterior trees using the program PACT v0.9.5<sup>32</sup> to compute the following summary statistics: a) genealogical diversity<sup>14</sup>, measuring the average time it takes for two randomly chosen contemporaneous lineages to coalesce, b) time to the most recent common ancestor (TMRCA)<sup>14</sup>, measuring the average time it takes for all contemporaneous lineages to find a common ancestor, c) genealogical  $F_{ST}$ , measuring the degree of population structure in contemporaneous lineages calculated as  $F_{ST} = (\pi_b - \pi_w)/\pi_b$ , where  $\pi_w$  is genealogical diversity between randomly sampled lineages from the same geographic region and  $\pi_b$  is genealogical diversity between randomly sampled lineages from different geographic regions, d) persistence, measuring the average number of years for a tip to leave its sampled location, walking backwards up the phylogeny, e) migration rate, measuring the average number of migration events over the phylogeny divided by total tree length to give migration events per lineage per year, f) trunk location through time<sup>4</sup>, measuring the posterior distribution across sampled phylogenies of the trunk geographic



state, where the trunk is defined as all branches ancestral to viruses sampled within 1 year of the most recent sample, g) region-specific ancestral geographic history, measuring the distribution of geographic locations of tips belonging to a particular region traced backwards in time through the phylogeny averaged across sampled phylogenies. Statistics (a), (b), (c), (f), and (g) were calculated across 0.1 year genealogical windows. These procedures gave an estimate of credible intervals for inferred ancestral locations across posterior phylogeographic reconstructions.

### Code and data availability

Sequence data has been deposited with the Influenza Research Database<sup>10</sup> and accession numbers provided as Source Data. The entire bioinformatic pipeline, including data subsampling, preparing XML files for BEAST, setting up PACT analyses and rendering figures is available at <https://github.com/blab/global-migration>. Analysis and data files are archived on the Dryad Digital Repository under DOI 10.5061/dryad.pc641.

### Surveillance, travel and age-structure data

We investigated epidemic size and frequency using virological isolation data between 2000 and 2012 collected by the WHO Collaborating Centre for Reference and Research on Influenza at the Victorian Infectious Diseases Reference Laboratory (VIDRL), Melbourne, Australia and the Centers for Disease Control and Prevention, Atlanta, USA (Extended Data Fig. 5f–i). These isolations were categorized by date of sampling and by virus type: H3N2, H1N1, Vic, or Yam. The data from VIDRL also contained information on patient age. The age structure of incidence was estimated by constructing a distribution of age of infection from individuals >5 y (owing to the overrepresentation of <5 year old patients for all subtypes) (Extended Data Fig. 5b–d). Median age of infection was 30 y (H3N2), 20 y (H1N1) and 16 y (B) and mean age of infection was 33.9 y (H3N2), 23.1 y (H1N1) and 23.2 y (B). Median age of infection was significantly different for H3N2 vs H1N1 ( $P = 4.6 \times 10^{-29}$ , Mann-Whitney  $U$  test), H3N2 vs B ( $P = 1.2 \times 10^{-62}$ ) and H1N1 vs B ( $P = 0.041$ ). The patient age data from VIDRL were potentially biased by testing strategy and the generally higher severity of H3N2 virus infections. Children and working age adults were more likely to be tested than the elderly but the greater severity of H3N2 virus infections might spread and flatten the patient age distribution. For this reason we additionally tested excluding individuals >65 y and recalculating summary statistics, finding median ages of infection of 27 y (H3N2), 19 y (H1N1) and 15 y (B) and mean age of infection as 28.0 y (H3N2), 22.2 y (H1N1) and 20.3 y (B). We classified children as 0–15 years and adults as 16 years and older, and estimated proportion of childhood infections as 30% (H3N2), 52% (H1N1) and 60% (B). There are potentially other biases specific to individual sentinel physicians and hospitals that could affect sample collection. However, the estimate derived from the VIDRL data that ~60% of influenza B virus infections occur in children is consistent with other estimates (reviewed in Glezen et al.<sup>8</sup>). Other studies similarly corroborate the estimates of lower age of infection for H1N1 viruses as compared to H3N2<sup>33,34</sup>.

Additionally, we analyzed the distribution of ages of ~102.5 million air passengers traveling through London Heathrow and London Gatwick airports in 2011 (Extended Data Fig. 5E)

reported by Civil Aviation Authority of the UK<sup>35</sup>. Assuming that children of ages 0 to 15 make up 17% of the UK population (Office of National Statistics), this distribution suggests that children engage in air travel at 19% the rate of adults.

For the modeling described below, we estimated age-structured contact rates following the empirical mixing data provided by Mossong et al.<sup>36</sup>. These contact matrices were previously validated in modeling pertussis epidemiology<sup>37</sup>. We simplified the Mossong et al. mixing matrices to record child-to-child contacts, child-to-adult contacts, adult-to-child contacts and adult-to-adult contacts, where children were defined to be 0 to 15 and adults to be 16 or over. This resulted in the following mixing matrix

$$\alpha = \begin{pmatrix} 1.0 & 0.21 \\ 0.21 & 0.26 \end{pmatrix},$$

where rates are relative to child-to-child contact rates.

### Epidemiological modeling

An individual-based model of influenza evolution and epidemiology was constructed following methods presented in Bedford et al.<sup>38</sup>. The model used here is identical to Bedford et al. except where specified below. The present implementation used a linear-strain space<sup>39,40</sup>, in which virus phenotype is represented by a continuous variable and cross-immunity between viruses is a function of distance between viruses in this space. We parameterized the model to compare scenarios of age-structured mixing between regions and to compare viruses with different rates of antigenic drift.

The model was simulated for 120 years with daily time steps and the first 100 years discarded to allow equilibrium to be reached. We modeled a metapopulation with individuals equally divided into three regions (North, Tropics, South). Individual's ages were tracked throughout the simulation and those less than 16 years old were classified as children and those 16 or older were classified as adults. Transmission occurred by mass action, with transmission rates modified by regional compartment and by age compartment. Thus, for example, the force of infection into children in the Tropics followed

$$\lambda_{ct} = \sum_{i \in (a,c)} \beta_t \alpha_{ic} I_{it} \frac{S_{ct}}{N_t} + \sum_{i \in (a,c)} \sum_{j \in (n,s)} \beta_t \alpha_{ic} m_i I_{ij} \frac{S_{ct}}{N_t},$$

where  $\beta_j$  is the seasonally forced contact rate in region  $j$ ,  $\alpha_{ac}$  represents adult-to-child transmission,  $m_i$  represents between-region transmission in age class  $i$ ,  $I_{ij}$  represents the number of infecteds in age class  $i$  in region  $j$ ,  $S_{ij}$  represents the number of susceptibles in age class  $i$  in region  $j$ , and  $N_j$  represents the total number of hosts in region  $j$ . The northern and southern regions were seasonally forced in opposite phase with a sinusoidal function following  $\varepsilon$ , while the tropics had no seasonal forcing.

Each virus possessed a one-dimensional antigenic phenotype  $\phi_v$  and after recovery a host 'remembered' its infecting phenotype. For each contact event, the Euclidean distance from

infecting phenotype  $\phi_v$  was calculated to each of the phenotypes in the host immune history  $\phi_{h_1}, \dots, \phi_{h_n}$ . Here, one unit of antigenic distance was designed to roughly correspond to a twofold dilution of antiserum in a hemagglutination inhibition (HI) assay<sup>41</sup>. The probability  $\rho$  that infection occurred after exposure was proportional to the distance  $d$  to the closest phenotype in the host immune history, following  $\rho = \min\{d s, 1\}$ . Each day there was a chance  $\mu$  that an infection mutates to a new phenotype. This mutation rate represents a phenotypic rate, rather than genetic mutation rate, and can be thought of as arising from multiple genetic sources. When a mutation occurred, the virus's phenotype was moved either left or right randomly and mutation size sampled from an exponential distribution with mean step size  $\delta_{\text{avg}}$ . Epidemiological parameters for the baseline epidemiological scenario with notation following Bedford et al.<sup>38</sup> were:

- Base transmission rate  $\beta = 0.88$  per day
- Duration of infection  $1/\nu = 5$  days
- Birth/death rate  $= 1/50$  years
- Total population size  $N = 45$  million
- Seasonal forcing in north and south  $\varepsilon = 0.15$
- Antigenic scaling  $s = 0.07$
- Antigenic mutation rate  $\mu = 0.5$  to  $6.5 \times 10^{-4}$  per day
- Average mutation size  $\delta_{\text{avg}} = 0.3$  units
- Child-to-child transmission  $\alpha_{cc} = 1.00$
- Child-to-adult transmission  $\alpha_{ca} = 0.21$
- Adult-to-child transmission  $\alpha_{ac} = 0.21$
- Adult-to-adult transmission  $\alpha_{aa} = 0.26$
- Child between-region transmission  $m_c = 0.0020$
- Adult between-region transmission  $m_a = 0.0020$

In the model with age-stratified mixing with host movement derived from air travel passenger age data, child between-region transmission  $m_c$  was 0.0011 and adult between-region transmission  $m_a$  was 0.0060.

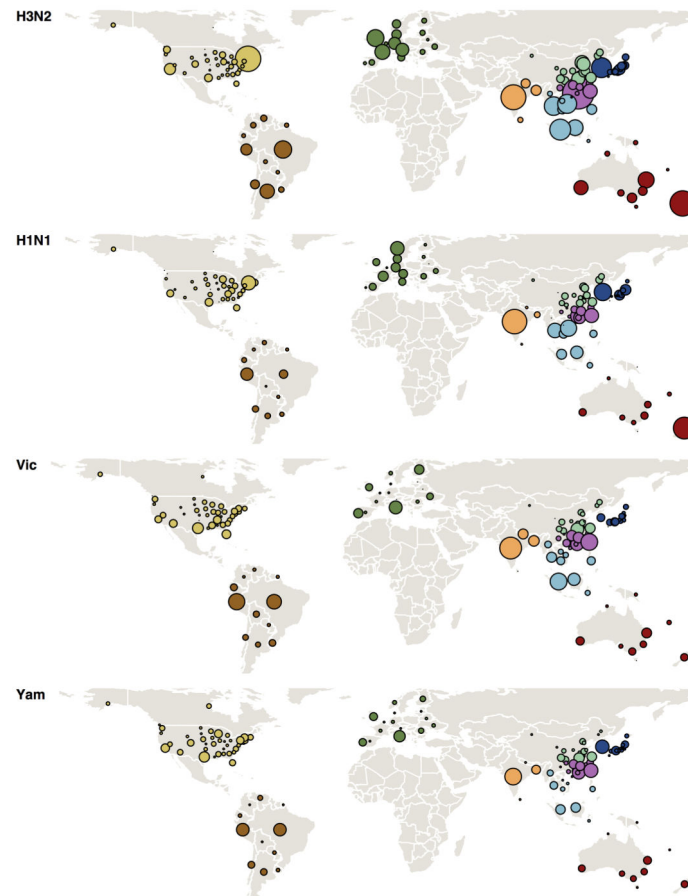
In the course of the simulation, the underlying infection history of who infects whom was recorded and output as a complete infection tree. Without ample within-host diversity owing to chronic infection, the complete infection tree also generated a fully observed phylogenetic tree. Examining geographic location across the phylogenetic tree allowed us to directly calculate migration rate as total migration events observed (transitions from one region to another) divided by total opportunity (tree length).

The simulation was parameterized to model H3-like, H1-like and B-like behavior (Extended Data Fig. 7) by modulating antigenic mutation rate  $\mu$  in the primary analysis (Fig. 3) or transmission rate  $\beta$  as a secondary analysis (Extended Data Fig. 8b). Values for  $\mu$  and  $\beta$  were

chosen based on observed attack rate, proportion of childhood infections and antigenic drift rate.

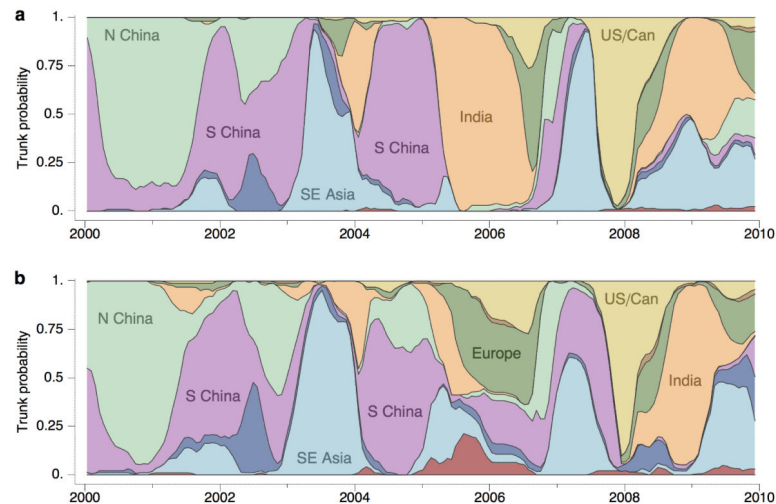
Source code for the simulation is available at <https://github.com/trvr/antigen/tree/global-migration> and parameter and results files are available at <https://github.com/blab/global-migration/tree/master/model>.

## Extended Data

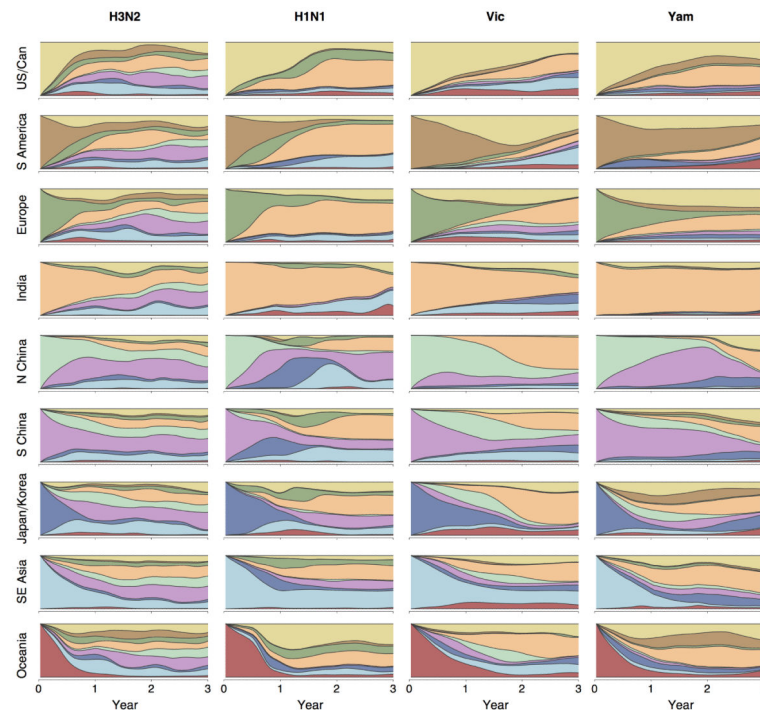


**Extended Data Figure 1. Spatial distribution of 4006 H3N2, 2144 H1N1, 1999 Vic and 1455 Yam samples**

Circle area is proportional to the number of sequenced viruses originating from a location. Color indicates assignment to one of 9 geographic regions.



**Extended Data Figure 2. Inferred location of the trunk of H3N2 tree through time in the primary dataset (a) and in a smaller secondary dataset (b)**  
Colored width at each time point indicates the posterior support for viruses from a particular geographic location comprising the trunk of the phylogenetic tree. Colors correspond to colored circles in persistence insets in Figure 1. The secondary datasets consist of 1391 H3N2 viruses, 1372 H1N1 viruses, 1394 Vic viruses and 1240 Yam viruses.

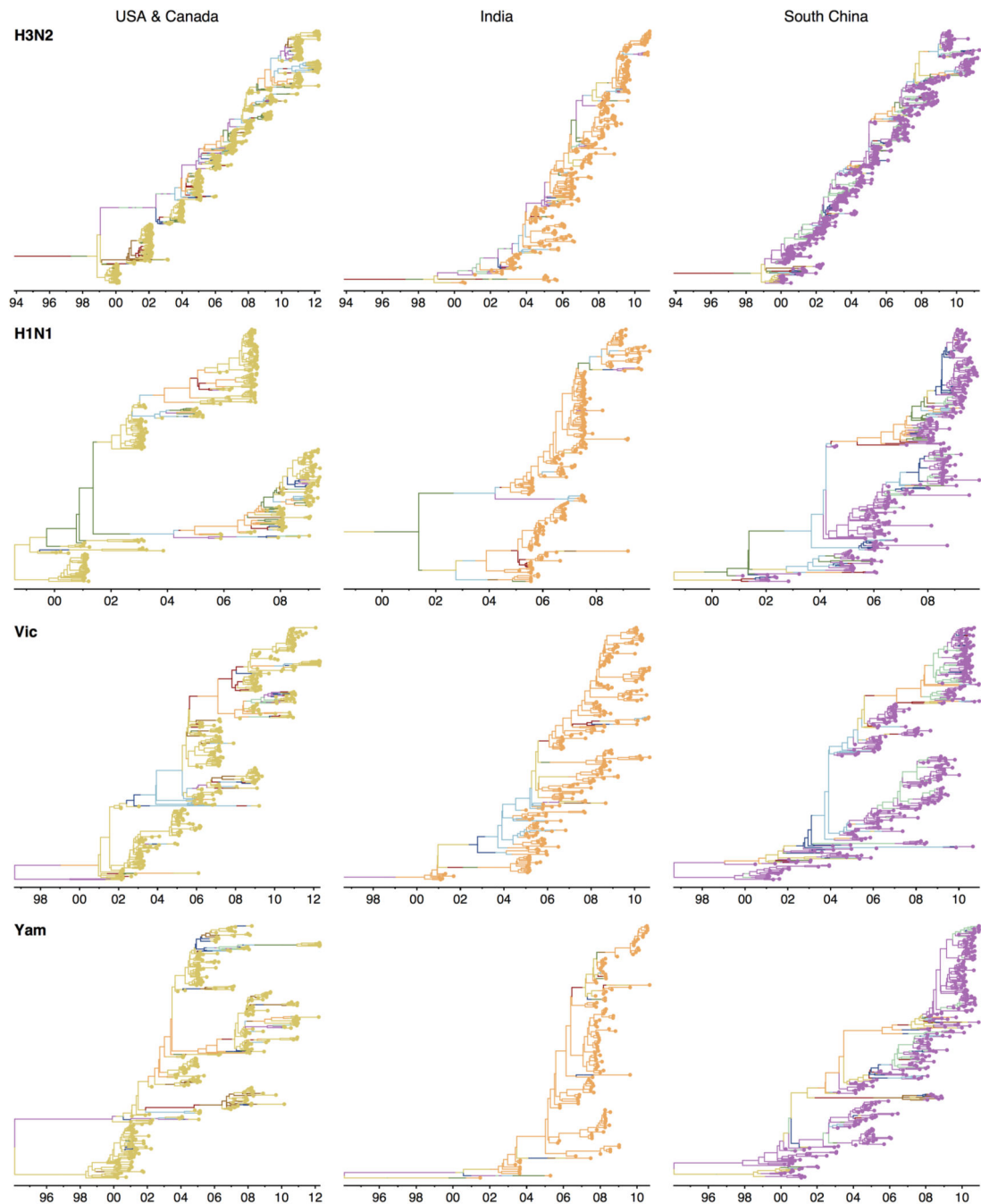


**Extended Data Figure 3. Average inferred geographic history of region-specific samples for H3N2, former seasonal H1N1, Vic and Yam viruses from 2000 to 2012**

In each panel, phylogeny tips belonging to a particular region were collected and their phylogeographic histories traced backwards in time averaging across the phylogenetic tree to combine all viruses within each region. The *x*-axis shows number of years backward in

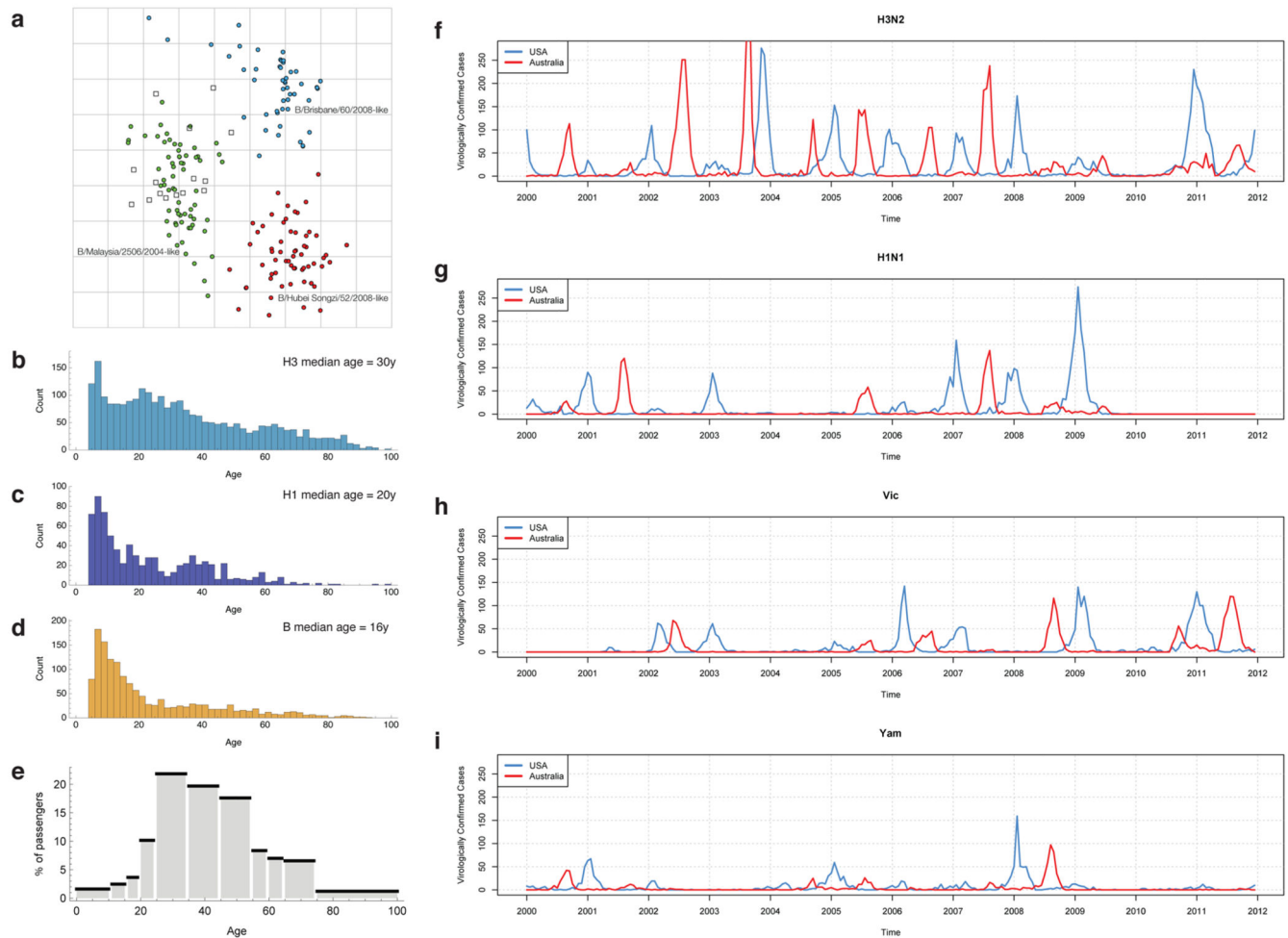
time from phylogeny tips from a particular region and the y-axis shows the geographic make up as stacked histogram of the ancestors of these tips, where region color-coding corresponds to the legend in Figure 1. For example, the top left panel shows the ancestry of USA and Canadian H3N2 viruses. At  $x = 0$ , all of these viruses are still in the USA or Canada and so an unbroken yellow band takes up the entire y. However, at  $x = 1$  year, a number of different geographic regions appear on the y. This indicates that, 1 year back, ancestors of USA and Canadian viruses are primarily found in Southeast Asia, India and South China. The pattern in the top right panel shows that the ancestors of USA and Canadian Yam viruses more often remain in the USA or Canada with approximately 50% of ancestors remaining 1 year back. Each panel is constructed by averaging across region-specific tips within a tree, but also across sampled posterior trees.



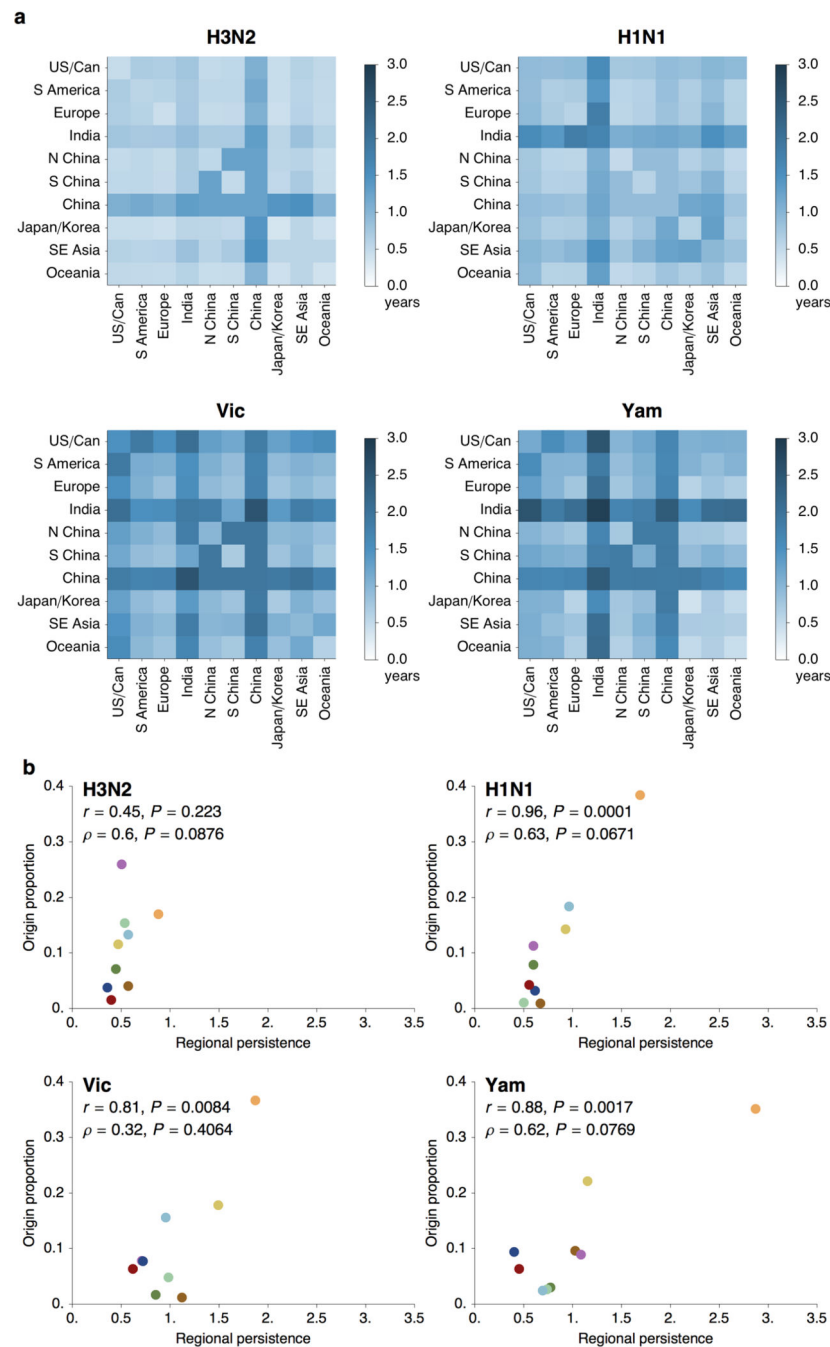


**Extended Data Figure 4. Maximum clade credibility (MCC) trees for region-specific samples from USA/Canada, India and South China for H3N2, H1N1, Vic and Yam viruses**

Each tree only contains viruses from a particular geographic region and thus tips are all a single color within a tree. Branch and trunk coloring have been retained from Figure 1 to highlight the inferred geographic ancestry of each lineage.



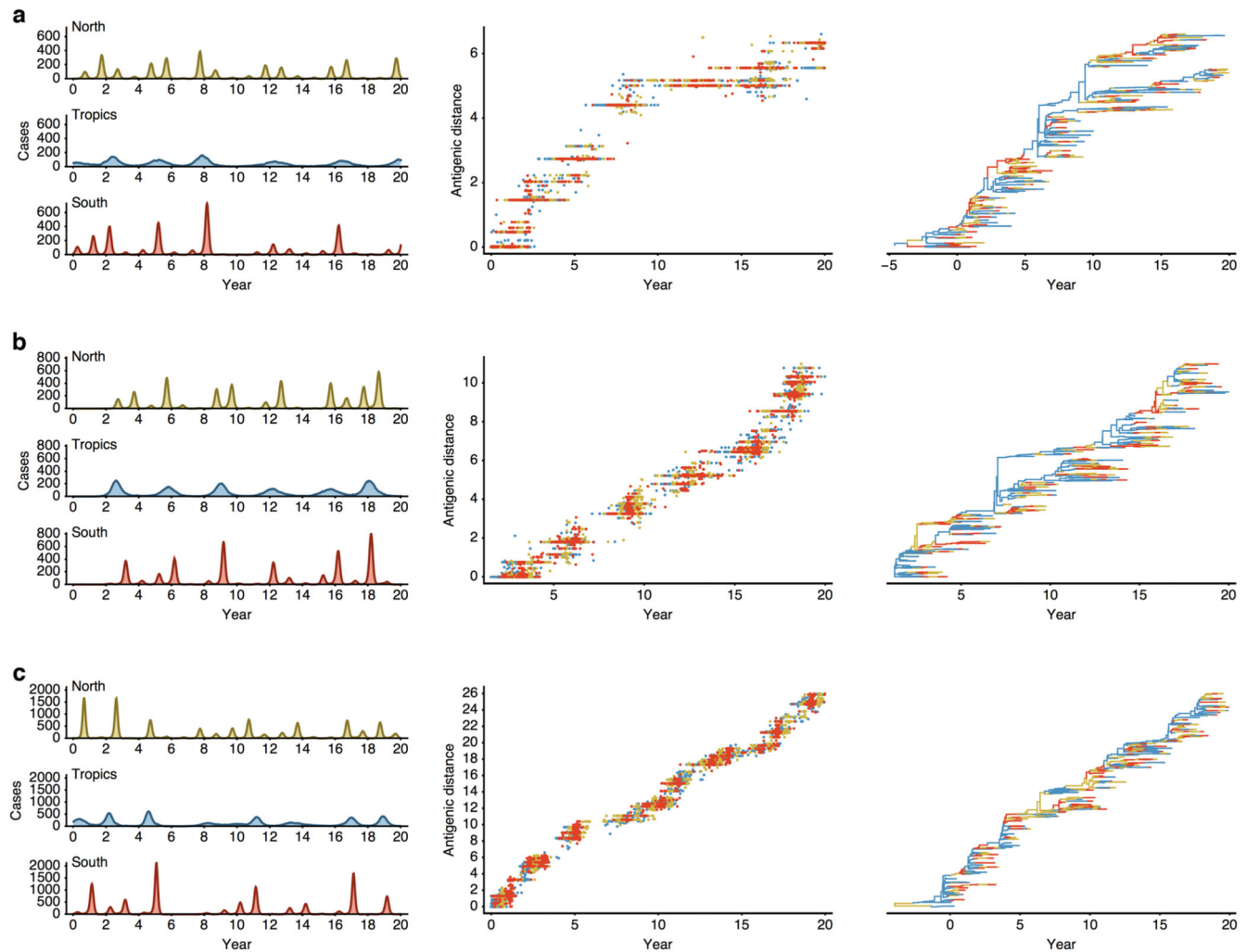
**Extended Data Figure 5. Antigenic map of Vic viruses primarily collected in 2008 (a), age distribution of infections for H3N2 (b), H1N1 (c) and B (d) in Australia 2000–2011, age distribution of ~102.5 million passengers at London Heathrow and London Gatwick airports during 2011 (e), timeseries of virological characterizations from 2000 to 2012 of viruses from the USA by US CDC and from Australia by VIDRL for H3N2 (f), H1N1 (g), Vic (h) and Yam (i)** In (a), the positions of strains (colored circles) and antisera (uncolored squares) are fit such that the distances between strains and antisera in the map represent the corresponding hemagglutination inhibition (HI) measurements with the least error following Smith et al.<sup>41</sup> using data on Vic viruses from the WHO Collaborating Centre for Reference and Research on Influenza at the Centers for Disease Control and Prevention, Atlanta, Georgia, USA. Strains are colored by antigenic cluster. Genetic clades corresponding to each antigenic cluster are marked with colored vertical bars in Fig 1c. The spacing between grid lines is one unit of antigenic distance corresponding to a twofold dilution of antiserum in the HI assay. In (f) to (i), virological characterizations are a surrogate for epidemiological activity that allow for accurate discrimination among H3N2, H1N1, Vic, and Yam viruses. These data generally reflect the relative magnitudes and frequencies of epidemics but in some cases will inflate magnitudes of very small epidemics due to preferential characterization of subtypes circulating at low levels.



**Extended Data Figure 6. Combined persistence estimates across pairs of regions for H3N2, H1N1, Vic and Yam (a) and Spearman correlation of a region's persistence vs the region's contribution to phylogenetic ancestry for H3N2, H1N1, Vic and Yam (b)**

In (a) and (b), persistence is measured as the average waiting time in years for a sample to leave its origin backwards in time in the phylogeny, with waiting time averaged across tips within a tree and across sampled posterior trees. In each panel of (a), the diagonal shows persistence within each of the 9 study regions and within the combined region of 'China', for which nodes in North China and in South China were considered to belong to a single region. The estimates along the diagonal are equivalent to the means shown in Figure 1. Off-

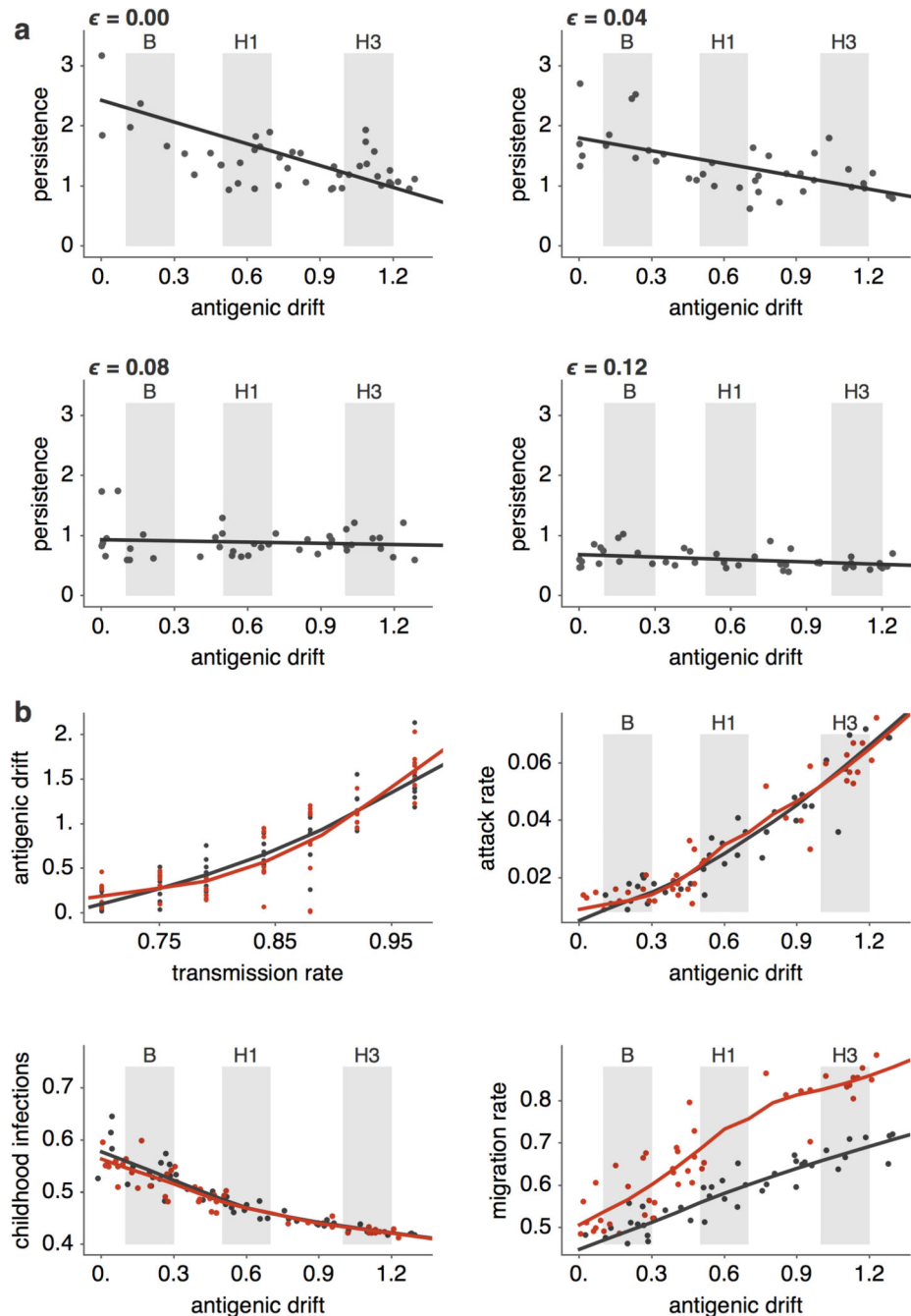
diagonal elements show persistence estimates for pairwise combinations of regions. For example, the off-diagonal for North and South China is exactly equivalent to the diagonal element for 'China' and the off diagonal for 'China' and India represents mean persistence when combining nodes from North China, South China and India. In (b), origin proportion is measured as the proportion of the time that a region is represented when tracing back 2 or more years from each tip in the phylogeny, averaged across tips within a tree and across sampled posterior trees. Spearman's  $\rho$  is not significant for any individual virus. However, the probability of observing 4 instances where each virus has a  $\rho$  of at least 0.32 is significant ( $P = 0.0017$ , bootstrap resampling test).



**Extended Data Figure 7. Simulation results for a model parameterized for slow antigenic drift (a), moderate antigenic drift (b), and fast antigenic drift (c)**

Colors represent geographic regions with tropics in blue, north in yellow and south in red. Region-specific incidence patterns are shown in terms of cases per 100,000 individuals per week, patterns of antigenic drift in terms of increasing antigenic distance (roughly proportional to  $\log_2$  HI units) over time and in the geographically labeled phylogeny. The parameterized antigenic mutation rate is 0.00015 antigenic mutations per infection per day

in (a), 0.00035 in (b) and 0.00055 in (c), while the realized antigenic drift rate is 0.29 antigenic units per year in (a), 0.58 in (b) and 1.19 in (c). Between-region mixing is 5.26 $\times$  faster in adults. Each panel shows output from a single simulation selected from the 112 shown in Figure 3, and is intended to show model behaviors over a range of parameters, not necessarily the behavior of particular viruses.



**Extended Data Figure 8. Simulation results showing relationship between antigenic drift and persistence as a function of seasonality (a) and simulation results showing the effects of modulating transmission rate  $\beta$  on model behavior (b)**

In (a), the seasonal forcing parameter  $\varepsilon$  follows  $\varepsilon = 0.00$  (no forcing),  $\varepsilon = 0.04$ ,  $\varepsilon = 0.08$  and  $\varepsilon = 0.12$  (moderate seasonal forcing). Points represent outcomes from a model in which adults travel between regions at  $5.26\times$  the rate of children. Solid black lines represent linear fits to the data. With 4 seasonality scenarios, 7 mutation rates and 8 replicates, there are 224 individual simulations shown. Persistence is measured as the average time in years taken for a tip to leave its region of origin going backwards in time, up the tree. In (b), transmission rate  $\beta$  in contacts per day is varied and compared to its effect on observed antigenic drift (in antigenic units per year), attack rate per year, proportion of childhood infections and migration rate between regions (in events per viral lineage per year). One antigenic unit is roughly equivalent to one  $\log_2$  HI unit. Black points represent outcomes from a model in which children and adults travel between regions at equal rates. Red points represent outcomes from a model in which adults travel between regions at  $5.26\times$  the rate of children. Solid black and red lines represent LOESS fits to the data. With 2 travel scenarios, 7 transmission rates and 8 replicates, there are 112 individual simulations shown.

**Extended Data Table 1**  
**Posterior mean estimates (and 95% highest posterior density intervals) across viruses for evolutionary and phylogeographic parameters**

Statistic	H3N2	H1N1	Vic	Yam
Total nucleotide rate <sup>*</sup>	5.0 (4.8–5.2)	4.4 (4.2–4.6)	2.7 (2.6–2.9)	2.8 (2.6–3.0)
Nonsynonymous rate <sup>*</sup>	2.2 (2.2–2.3)	1.9 (1.9–2.0)	1.0 (0.9–1.1)	1.0 (0.9–1.0)
Synonymous rate <sup>*</sup>	2.8 (2.7–2.9)	2.6 (2.5–2.7)	1.8 (1.8–1.9)	1.8 (1.8–1.9)
Antigenic drift rate <sup>†</sup>	1.01 (0.98–1.04)	0.62 (0.56–0.67)	0.42 (0.32–0.52)	0.32 (0.25–0.39)
Diversity <sup>‡</sup>	3.03	4.59	5.46	6.83
TMRCA <sup>§</sup>	3.89	4.53	5.22	7.62
$F_{ST}$ <sup>  </sup>	0.30	0.36	0.37	0.36
Persistence <sup>¶</sup>	0.50 (0.48–0.54)	0.79 (0.73–0.85)	1.07 (0.98–1.16)	1.03 (0.88–1.21)
Migration rate <sup>#</sup>	1.99 (1.85–2.10)	1.27 (1.18–1.37)	0.93 (0.86–1.02)	0.98 (0.83–1.14)

<sup>\*</sup> Evolutionary rates are measured in terms of  $10^{-3}$  substitutions per site per year.

<sup>†</sup> Antigenic drift rates are from Bedford et al.<sup>13</sup> table 2, and measures cartographic drift per year in terms of twofold dilution of antiserum in a hemagglutination inhibition (HI) assay.

<sup>‡</sup> Diversity of contemporaneous lineages is measured as average time in years for two randomly sampled lineages to share a common ancestor.

<sup>§</sup> Time to the most recent common ancestor (TMRCA) of contemporaneous lineages is measured as the average time in years for all lineages to find a common ancestor.

<sup>||</sup>  $F_{ST}$  compares diversity within regions to diversity between regions, so that  $F_{ST} = (\pi_b - \pi_w) / \pi_b$ .

<sup>¶</sup> Persistence is calculated as the average number of years for a tip to leave its sampled location, walking backwards up the phylogeny.

<sup>#</sup> Migration rate is calculated as migration events per lineage per year between any two regions.



**Extended Data Table 2**  
**Posterior mean estimates across viruses and datasets of**  
**regional persistence, migration rate and geographic**  
**population structure**

Statistic	Dataset	H3N2	H1N1	Vic	Yam
Persistence <sup>*</sup>	Primary <sup>§</sup>	0.51	0.79	1.07	1.03
Persistence <sup>*</sup>	Secondary <sup>//</sup>	0.53	0.75	1.16	1.11
Persistence <sup>*</sup>	Alternative <sup>¶</sup>	0.50	0.76	1.28	1.12
Migration rate <sup>†</sup>	Primary <sup>§</sup>	1.96	1.27	0.93	0.97
Migration rate <sup>†</sup>	Secondary <sup>//</sup>	1.89	1.33	0.86	0.90
Migration rate <sup>†</sup>	Alternative <sup>¶</sup>	2.00	1.32	0.78	0.89
$F_{ST}$ <sup>‡</sup>	Primary <sup>§</sup>	0.30	0.36	0.37	0.36
$F_{ST}$ <sup>‡</sup>	Secondary <sup>//</sup>	0.29	0.35	0.36	0.37
$F_{ST}$ <sup>‡</sup>	Alternative <sup>¶</sup>	0.29	0.34	0.36	0.35

<sup>\*</sup> Regional persistence is measured as the average waiting time in years for a sample to leave its origin backwards in time in the phylogeny.

<sup>†</sup> Migration rate is measured as migration events per lineage per year.

<sup>‡</sup>  $F_{ST}$  compares diversity within regions to diversity between regions, so that  $F_{ST} = (\pi_b - \pi_w) / \pi_b$ .

<sup>§</sup> The primary datasets consist of 4006 H3N2 viruses, 2144 H1N1 viruses, 1999 Vic viruses and 1455 Yam viruses.

<sup>//</sup> The secondary datasets consist of 1391 H3N2 viruses, 1372 H1N1 viruses, 1394 Vic viruses and 1240 Yam viruses.

<sup>¶</sup> The alternative datasets consist of 1967 H3N2 viruses, 1439 H1N1 viruses, 1756 Vic viruses and 1223 Yam viruses divided into 10 geographic regions (USA/Canada, South America, Europe, India, Japan/Korea, Southeast Asia, Oceania, China, Central America and Africa).

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

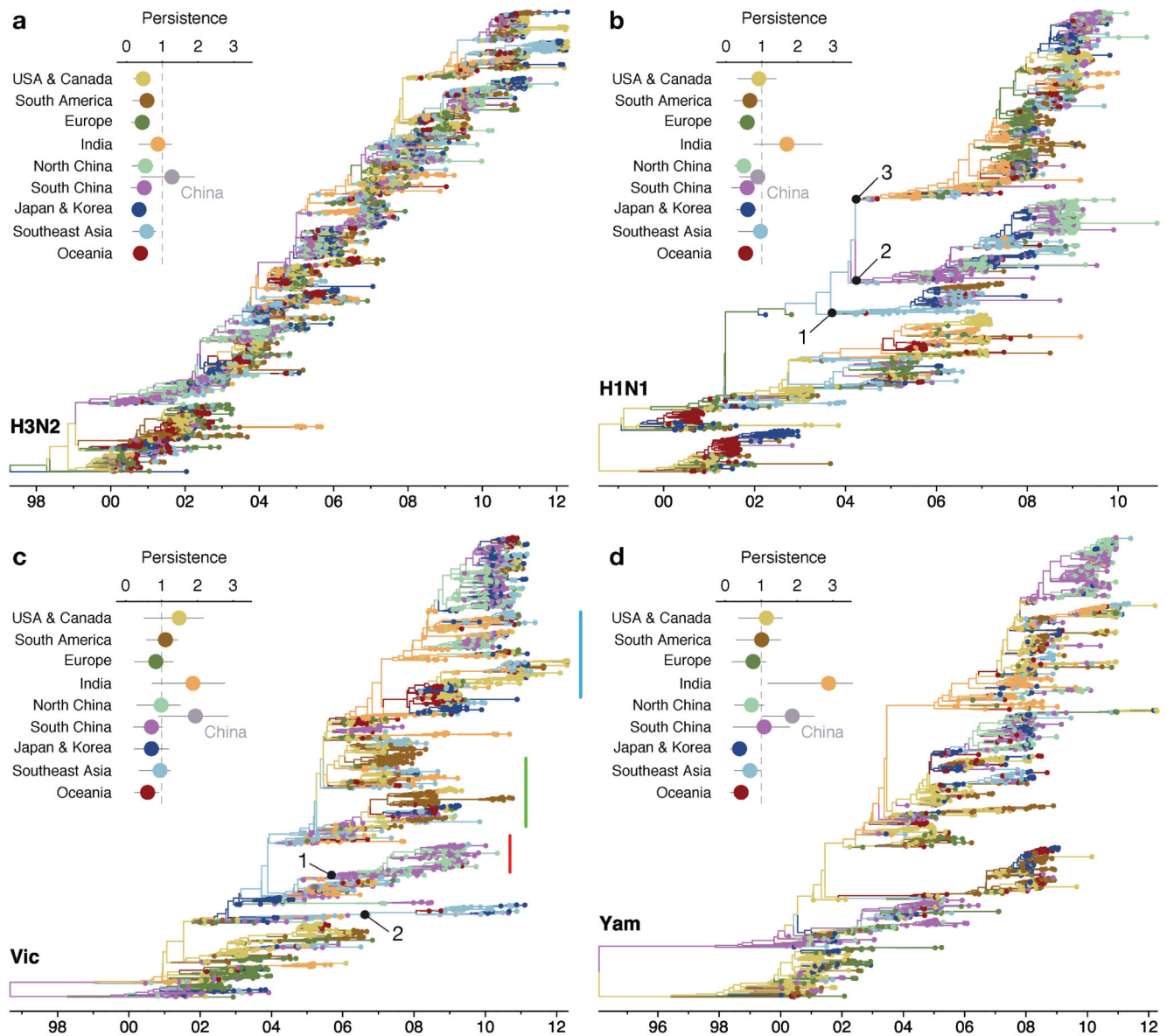
T.B. was supported by a Newton International Fellowship from the Royal Society and through NIH U54 GM111274. S.R. was supported by MRC (UK, Project MR/J008761/1), Wellcome Trust (UK, Project 093488/Z/10/Z), Fogarty International Centre (USA, R01 TW008246-01), DHS (USA, RAPIDD program), NIGMS (USA, MIDAS U01 GM110721-01) and NIHR (UK, Health Protection Research Unit funding). The Melbourne WHO Collaborating Centre for Reference and Research on Influenza was supported by the Australian Government Department of Health and thanks N. Komadina and Y.-M. Deng. The Atlanta WHO Collaborating Center for Surveillance, Epidemiology and Control of Influenza was supported by the U.S. Department of Health and Human Services. NIV thanks A.C. Mishra, M. Chawla-Sarkar, A.M. Abraham, D. Biswas, S. Shrikhande, AnuKumar B, and A. Jain. Influenza surveillance in India was expanded, in part, through US Cooperative Agreements (5U50C1024407 and U51IP000333) and by the Indian Council of Medical Research. M.A.S. was supported through NSF DMS 1264153 and NIH R01 AI 107034. Work of the WHO Collaborating Centre for Reference and Research on Influenza at the MRC National Institute for Medical Research was supported by U117512723. P.L., A.R. & M.A.S were supported by EU Seventh Framework Programme [FP7/2007-2013] under Grant Agreement no. 278433-PREDEMICS and ERC Grant agreement no. 260864. C.A.R. was supported by a University Research Fellowship from the Royal Society.

## References

1. Russell CA, et al. The global circulation of seasonal influenza A (H3N2) viruses. *Science*. 2008; 320:340–346. [PubMed: 18420927]

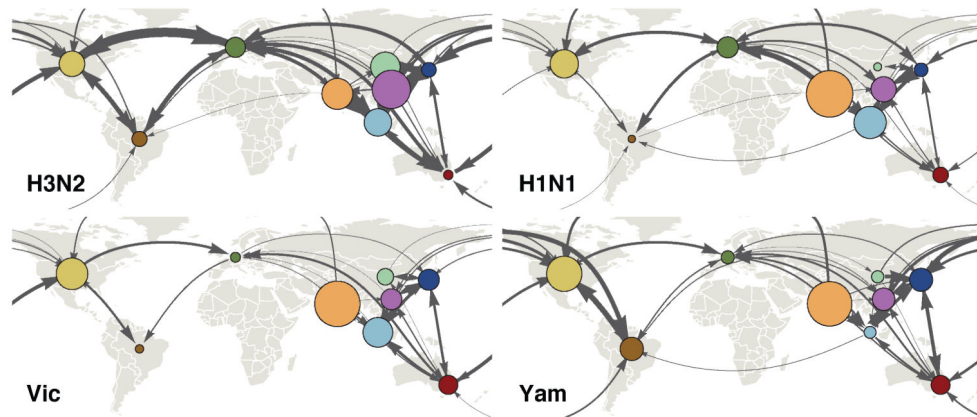
2. Lemey P, et al. Unifying viral genetics and human transportation data to predict the global transmission dynamics of human influenza H3N2. *PLoS Pathog.* 2014; 10:e1003932. [PubMed: 24586153]
3. Rambaut A, et al. The genomic and epidemiological dynamics of human influenza A virus. *Nature.* 2008; 453:615–619. [PubMed: 18418375]
4. Bedford T, Cobey S, Beerli P, Pascual M. Global migration dynamics underlie evolution and persistence of human influenza A (H3N2). *PLoS Pathog.* 2010; 6:e1000918. [PubMed: 20523898]
5. Chan J, Holmes A, Rabadan R. Network analysis of global influenza spread. *PLoS Comput Biol.* 2010; 6:e1001005. [PubMed: 21124942]
6. Nelson MI, Simonsen L, Viboud C, Miller MA, Holmes EC. Phylogenetic analysis reveals the global migration of seasonal influenza A viruses. *PLoS Pathog.* 2007; 3:1220–1228. [PubMed: 17941707]
7. Nelson MI, et al. Stochastic processes are key determinants of short-term evolution in influenza A virus. *PLoS Pathog.* 2006; 2:e125. [PubMed: 17140286]
8. Glezen PW, Schmier JK, Kuehn CM, Ryan KJ, Oxford J. The burden of influenza B: a structured literature review. *Am J Public Health.* 2013; 103:e43–51.
9. Thompson WW, et al. Influenza-associated hospitalizations in the United States. *Jama.* 2004; 292:1333–1340. [PubMed: 15367555]
10. Squires RB, et al. Influenza research database: an integrated bioinformatics resource for influenza research and surveillance. *Influenza Other Respir Viruses.* 2012; 6:404–416. [PubMed: 22260278]
11. Chen R, Holmes EC. The evolutionary dynamics of human influenza B virus. *J Mol Evol.* 2008; 66:655–663. [PubMed: 18504518]
12. Cox NJ, Bender CA. The molecular epidemiology of influenza viruses. *Seminars in Virology.* 1995; 6:359–370.
13. Bedford T, et al. Integrating influenza antigenic dynamics with molecular evolution. *eLife.* 2014; 3:e01914. [PubMed: 24497547]
14. Bedford T, Cobey S, Pascual M. Strength and tempo of selection revealed in viral gene genealogies. *BMC Evol Biol.* 2011; 11:220. [PubMed: 21787390]
15. Lemey P, Rambaut A, Drummond AJ, Suchard MA. Bayesian phylogeography finds its roots. *PLoS Comput Biol.* 2009; 5:e1000520. [PubMed: 19779555]
16. Fox JP, Hall CE, Cooney MK, Foy HM. Influenzavirus infections in Seattle families, 1975–1979. I. Study design, methods and the occurrence of infections by time and age. *Am J Epidemiol.* 1982; 116:212–227. [PubMed: 7114033]
17. Longini IM Jr, Koopman JS, Monto AS, Fox JP. Estimating household and community transmission parameters for influenza. *Am J Epidemiol.* 1982; 115:736–751. [PubMed: 7081204]
18. Viboud C, et al. Synchrony, waves, and spatial hierarchies in the spread of influenza. *Science.* 2006; 312:447–451. [PubMed: 16574822]
19. Recker M, Pybus OG, Nee S, Gupta S. The generation of influenza outbreaks by a network of host immune responses against a limited set of antigenic types. *Proc Natl Acad Sci U S A.* 2007; 104:7711–7716. [PubMed: 17460037]
20. Zinder D, Bedford T, Gupta S, Pascual M. The roles of competition and mutation in shaping antigenic and genetic diversity in influenza. *PLOS Pathog.* 2013; 9:e1003104. [PubMed: 23300455]
21. Nobusawa E, Sato K. Comparison of the mutation rates of human influenza A and B viruses. *J Virol.* 2006; 80:3675–3678. [PubMed: 16537638]
22. Neuzil KM, et al. Immunogenicity and reactogenicity of 1 versus 2 doses of trivalent inactivated influenza vaccine in vaccine-naïve 5–8-year-old children. *J Inf Dis.* 2006; 194:1032–1039. [PubMed: 16991077]
23. Matrosovich MN, et al. Probing of the receptor-binding sites of the H1 and H3 influenza A and influenza B virus hemagglutinins by synthetic and natural sialosides. *Virology.* 1993; 196:111–121. [PubMed: 8356788]
24. Hensley S, et al. Hemagglutinin receptor binding avidity drives influenza A virus antigenic drift. *Science.* 2009; 326:734–736. [PubMed: 19900932]

25. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* 2012; 29:1969–1973. [PubMed: 22367748]
26. Shapiro B, Rambaut A, Drummond A. Choosing appropriate substitution models for the phylogenetic analysis of protein-coding sequences. *Mol. Biol. Evol.* 2006; 23:7–9. [PubMed: 16177232]
27. URL <http://tree.bio.ed.ac.uk/software/pathogen/>
28. Shapiro B, et al. A Bayesian phylogenetic method to estimate unknown sequence ages. *Mol. Biol. Evol.* 2011; 28:879–887. [PubMed: 20889726]
29. Pagel M, Meade A, Barker D. Bayesian estimation of ancestral character states on phylogenies. *Syst. Biol.* 2004; 53:673–684. [PubMed: 15545248]
30. Lemey P, Minin VN, Bielejec F, Pond SLK, Suchard MA. A counting renaissance: combining stochastic mapping and empirical Bayes to quickly detect amino acid sites under positive selection. *Bioinformatics.* 2012; 28:3248–3256. [PubMed: 23064000]
31. Edwards CJ, et al. Ancient hybridization and an Irish origin for the modern polar bear matriline. *Curr. Biol.* 2011; 21:1251–1258. [PubMed: 21737280]
32. URL <https://github.com/trvr/PACT>
33. Kelly H, Grant K, Williams S, Fielding J, Smith D. Epidemiological characteristics of pandemic influenza H1N1 2009 and seasonal influenza infection. *Med. J. Aust.* 2009; 191:146–149. [PubMed: 19645642]
34. Khiabani H, Farrell G, St George K, Rabadan R. Differences in patient age distribution between influenza A subtypes. *PLOS ONE.* 2009; 4:e6832. [PubMed: 19718262]
35. URL <http://www.caa.co.uk/docs/81/2011CAAPaxSurveyReport.pdf>
36. Mossong J, et al. Social contacts and mixing patterns relevant to the spread of infectious diseases. *PLOS Med.* 2008; 5:e74. [PubMed: 18366252]
37. Rohani P, Zhong X, King AA. Contact network structure explains the changing epidemiology of pertussis. *Science.* 2010; 330:982–985. [PubMed: 21071671]
38. Bedford T, Rambaut A, Pascual M. Canalization of the evolutionary trajectory of the human influenza virus. *BMC Biol.* 2012; 10:38. [PubMed: 22546494]
39. Gog JR, Grenfell BT. Dynamics and selection of many-strain pathogens. *Proc. Natl. Acad. Sci. USA.* 2002; 99:17209–17214. [PubMed: 12481034]
40. Lin J, Andreasen V, Casagrandi R, A Levin S. Traveling waves in a model of influenza A drift. *J. Theor. Biol.* 2003; 222:437–445. [PubMed: 12781742]
41. Smith DJ, et al. Mapping the antigenic and genetic evolution of influenza virus. *Science.* 2004; 305:371–376. [PubMed: 15218094]



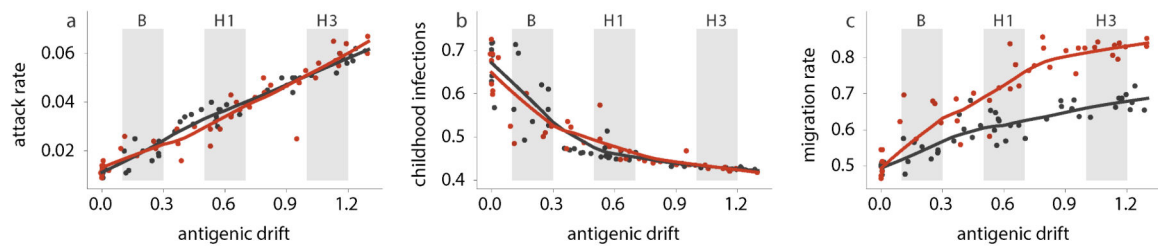
**Figure 1. Maximum clade credibility trees for primary datasets of 4006 H3N2 viruses (a), 2144 H1N1 viruses (b), 1999 Vic viruses (c) and 1455 Yam viruses (d)**

Branch tips are colored by geographic region of virus collection; internal branches are colored by geographic region as inferred by Bayesian phylogeographic methods (region colors in persistence insets). In b) nodes 1-3 indicate co-circulating clades that diverged in 2004. In c), nodes 1 and 2 indicate divergent clades of viruses from Asia, colored vertical bars indicate antigenic variants shown in Extended Data Figure 5a (green: B/Malaysia/2506/2004-like, red: B/Hubei Songzi/52/2008-like, other post-2008 viruses: B/Brisbane/60/2008-like). The inset to the top left of each tree shows duration of region-specific persistence measured as the waiting time in years for a virus to leave its geographic region of origin. Circles represent mean persistence across sampled viruses, while lines show the inter-quartile range of persistence across sampled viruses. Region "China", shows the combined persistence estimate for North China and South China together.



**Figure 2. Estimates of mean pairwise virus migration rate**

Line thickness between regions indicates average number of migration events per lineage per year. Arrowhead size indicates the strength of directionality of migration. For clarity, only arrows corresponding to migration rates greater than 0.25 events per lineage per year are shown. Circle area indicates the global proportion of ancestry deriving from each region.



**Figure 3. Relationship of antigenic drift to incidence (a), proportion of childhood infections (b), and geographic migration rate (c), in a multi-strain multi-region model of influenza transmission** Black points represent outcomes from a model in which children and adults travel between regions at equal rates. Red points represent outcomes from a model in which adults travel between regions at 5.26 $\times$  the rate of children (Extended Data Fig. 5e). Solid black and red lines represent LOESS fits to the data. With 2 travel scenarios, 7 mutation rates and 8 replicates, there are 112 individual stochastic simulations (Extended Data Fig. 7). Antigenic drift was measured in cartographic units<sup>13</sup> per year (see Methods). In a) attack rate was measured as proportion of the total population infected yearly. In c) migration rate was measured in terms of migration events per lineage per year.